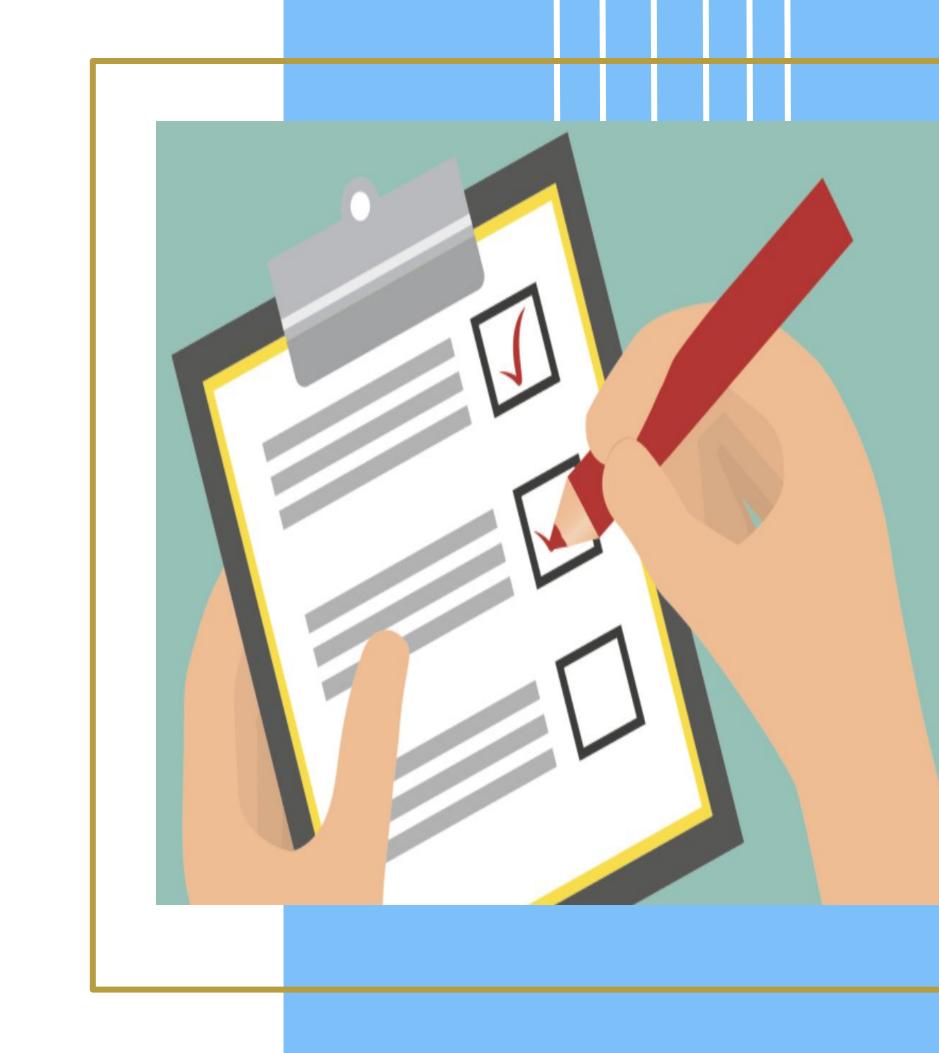
Inventa Analytics: criando vantagem competitiva de negócio com dados + IA

**Lucas Fonseca** 

#### Agenda e Pontos-Chave

- Apresentação pessoal
- Escopo da conversa
- Fundação de Dados
- Revisão Conceitual
- Inventa Analytics
- Código
- Limitações e Próx. Passos



#### Prazer, Lucas

- De Campos Novos, Cunha SP
- Formado em Análise e

  Desenvolvimento de Sistemas

  pela FATEC-SP em 2016
- Especialista em Inteligência e gestão de dados pela Escola Politécnica da USP (2018)
- Cursando Gestão Empresarial na EAESP FGV (2025)
- Há 12 anos no mercado de tecnologia. Experiência em Engenharia de SW e Dados, Analytics Engineering, Machine Learning Ops e liderança
- Membro do AGI Club (by Ifood)



Lavandário • Cunha-SP









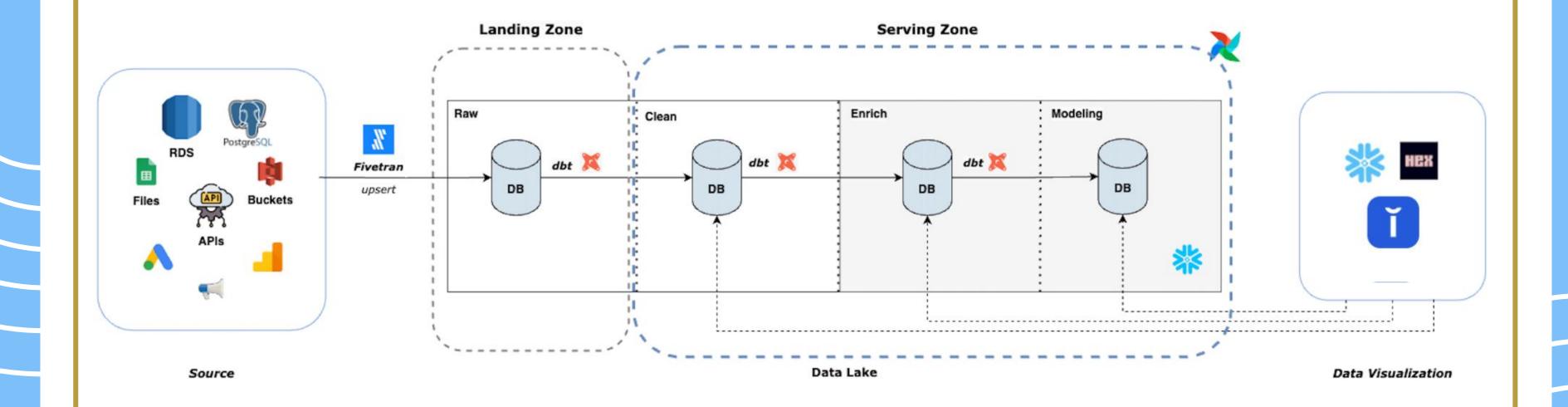
#### Escopo da conversa

- Conversaremos sobre como estamos construindo o Inventa Analytics, uma solução conversacional baseada em Dados+LLM para self-service em dados na Inventa (B2x Holding).
- Engenharia de dados como fundação para produtos de Dados e IA.
- Arquitetura modular de agentes para interpretação de intenções, geração queries e entrega de insights confiáveis.
- Tech Stack do projeto: Snowflake, OpenAI, FastAPI.
- Aprendizados em produtos inovadores, equilibrando confiabilidade e experimentação.
- Experiência do usuário (UX + AI) como parte central da arquitetura.
- ❖ Do PoC ao MVP até a evolução para produto em produção.

Engenharia de dados como fundação para produtos de Dados e IA

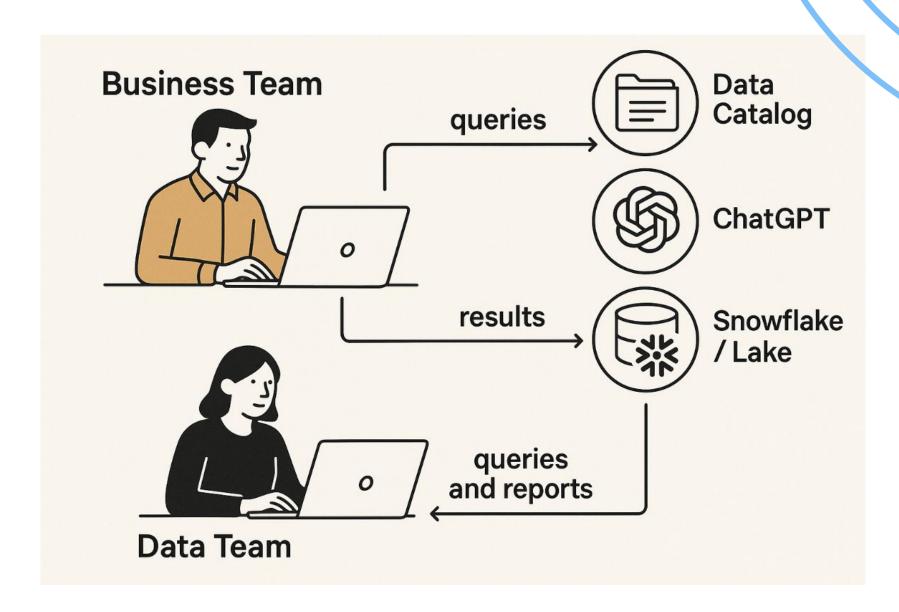


#### **Artemis Data Platform**



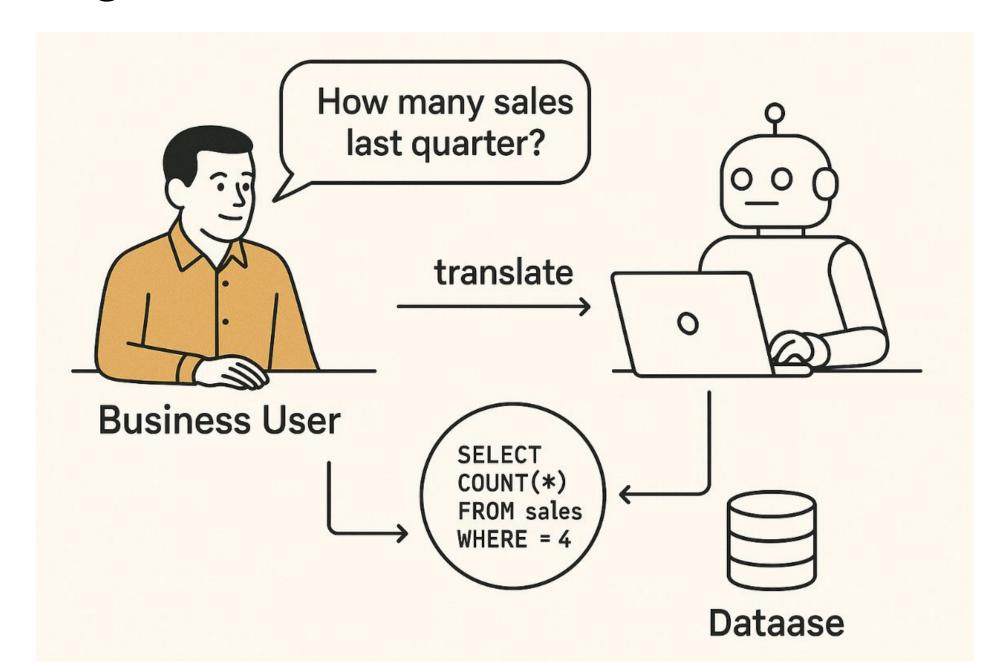
#### Problema

- Time de negócios não possui fluência em SQL para consultar os dados
- Time de Dados se tornou gargalo para a empresa
  - Manutenção da Plataforma
  - Criação de Produtos de Dados
  - Consultas de Dados

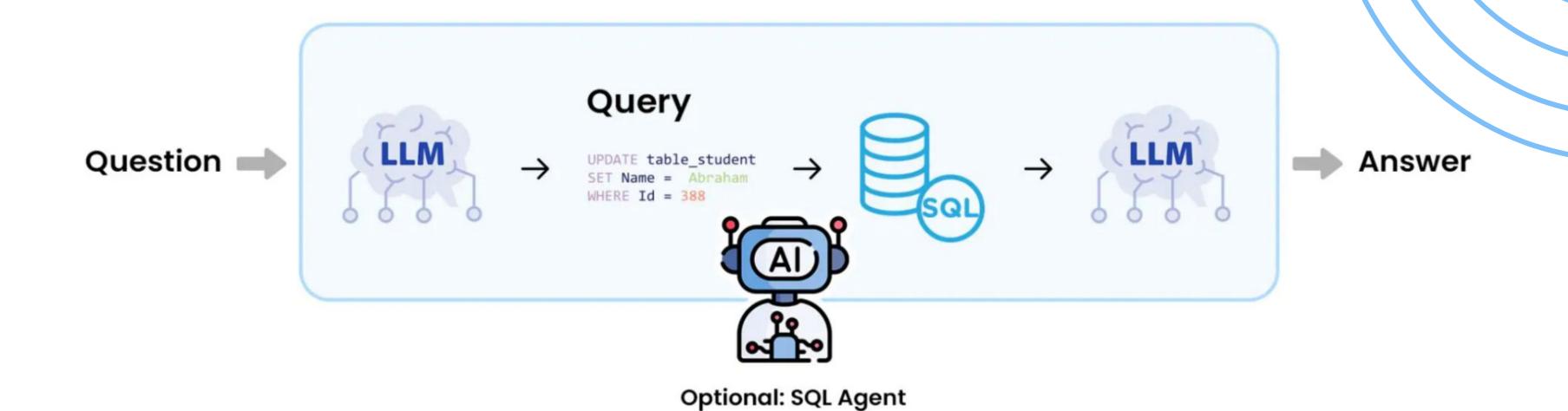


#### Problema

Ponto central: Como poderíamos democratizar os dados da Inventa através de tecnologia?



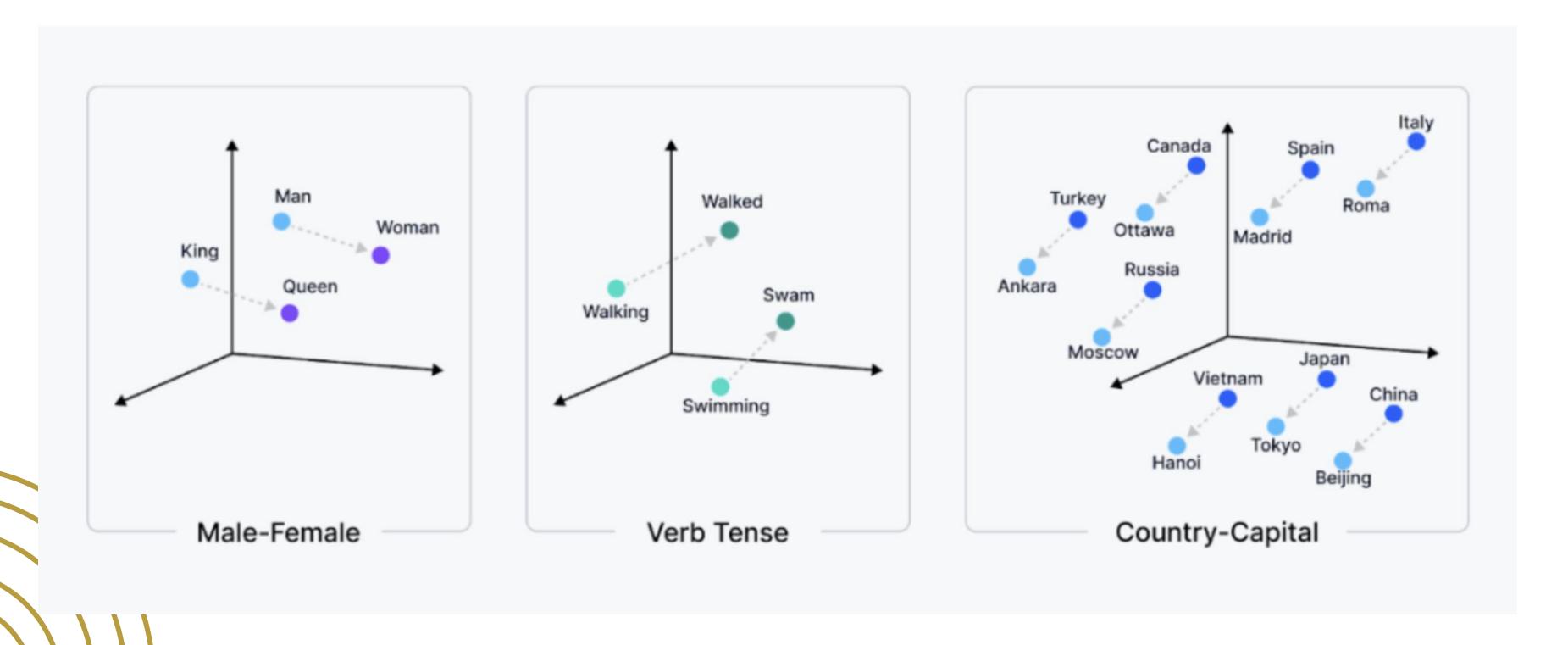
### Text-to-Sql/NL2SQL



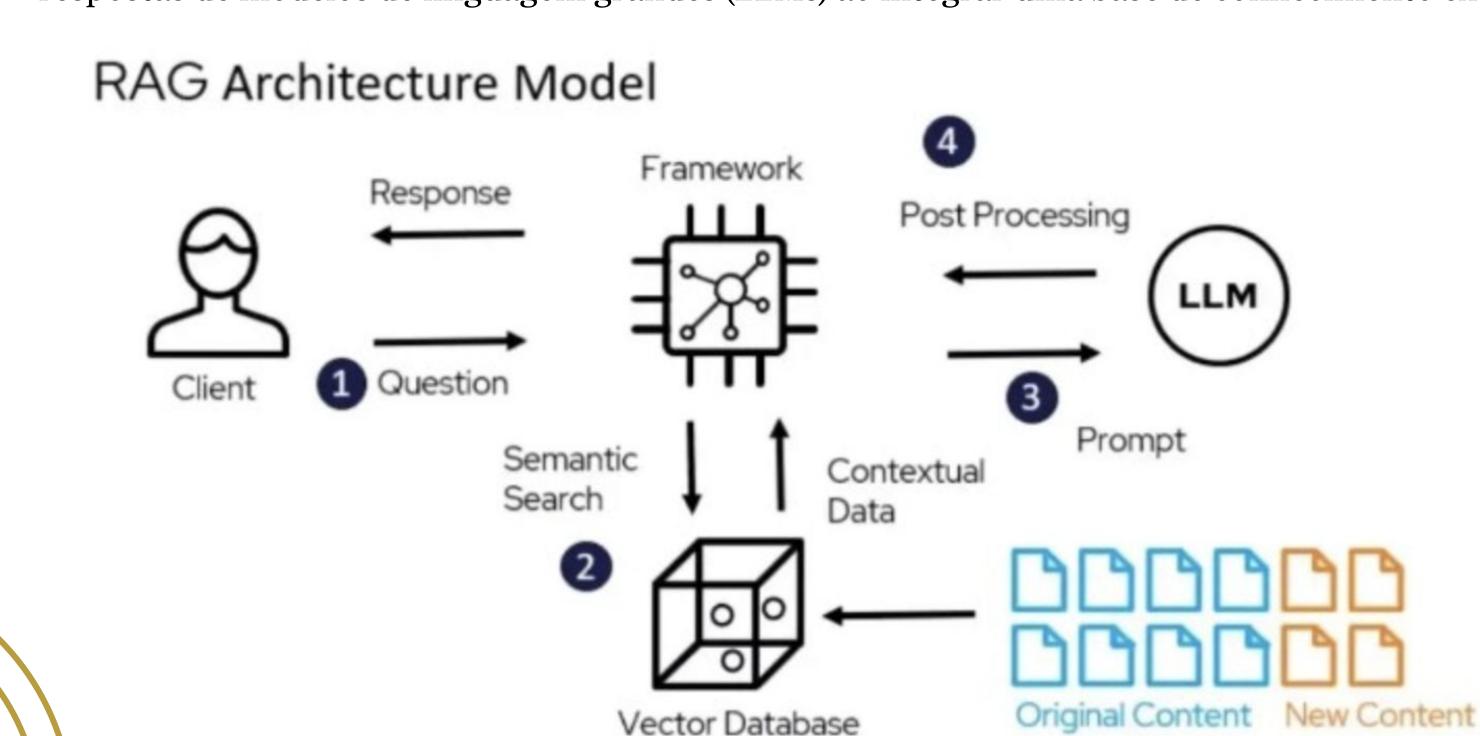
### Algumas definições conceituais

- **Text-to-sql:** É uma tecnologia capaz de converter perguntas em linguagem natural para comandos SQL.
- ♦ Modelos fundacionais: os modelos de inteligência artificial de uso geral, treinados com enormes volumes de dados (como textos, imagens e códigos). Ex.: GPT-5 (OpenAI), Gemini (Google), Claude (Anthropic), LLaMA (Meta), Mistral, etc.
- **Agentes:** Agente (em Inteligência Artificial) é uma entidade capaz de perceber o ambiente, **tomar decisões** e agir para atingir um **objetivo** .

Embedding: É a representação numérica de um dado textual (ou imagem, áudio, etc.) em um espaço matemático. Ex: carro e automóvel são embeddings parecidos.



\* RAG (Retrieval Augmented Generation): RAG (Retrieval-Augmented Generation, ou Geração Aumentada por Recuperação) é uma técnica de IA que melhora a precisão e relevância das respostas de modelos de linguagem grandes (LLMs) ao integrar uma base de conhecimento externa.



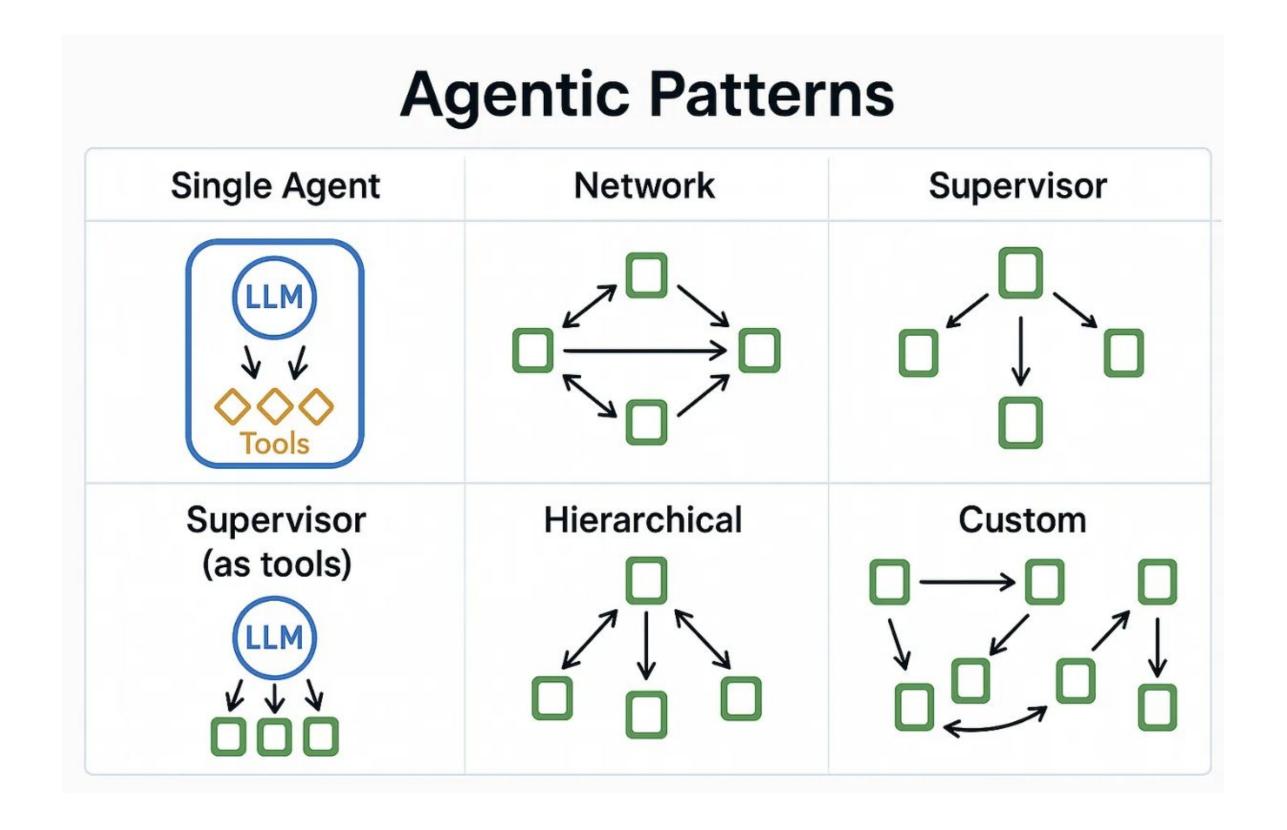
#### **Arquitetura Agentic:**

- > agentes são autônomos e decidem quando e com quem interagir
- > geram suas próprias mensagens
- > e avaliam respostas para continuar o raciocínio

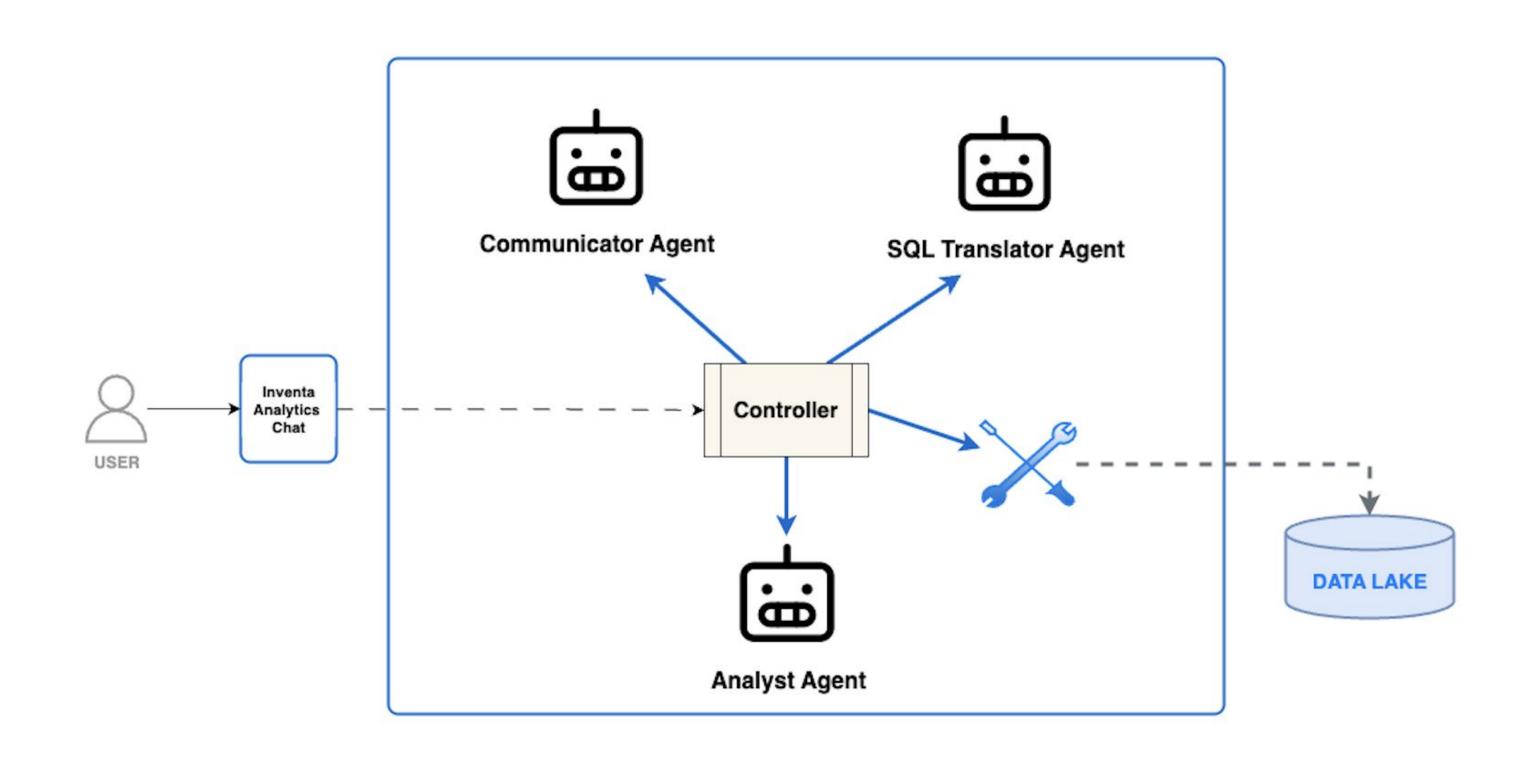
#### Arquitetura baseada em Agentes:

- > o controle de fluxo é determinístico
- > um controlador/router define a priori quem fala com quem
- > cada agente ainda possui autonomia local (decisão sobre o conteúdo da resposta), mas não têm liberdade de escolha sobre ações ou interações

## Arquitetura Agentic



### Arquitetura Baseada em Agentes



#### **Inventa Analytics**

#### Bem-vindo Lucas Mendes Mota da Fonseca



Qual o total de vendas por mes nesse ano corrente?



**Explicação da query:** A consulta retorna o GMV total por mês para o ano de 2025, considerando apenas pedidos válidos e não subsidiados. O agrupamento é feito pelo ano e mês da data de criação do pedido, e o resultado é ordenado cronologicamente.

ANO\_MES GMV\_TOTAL

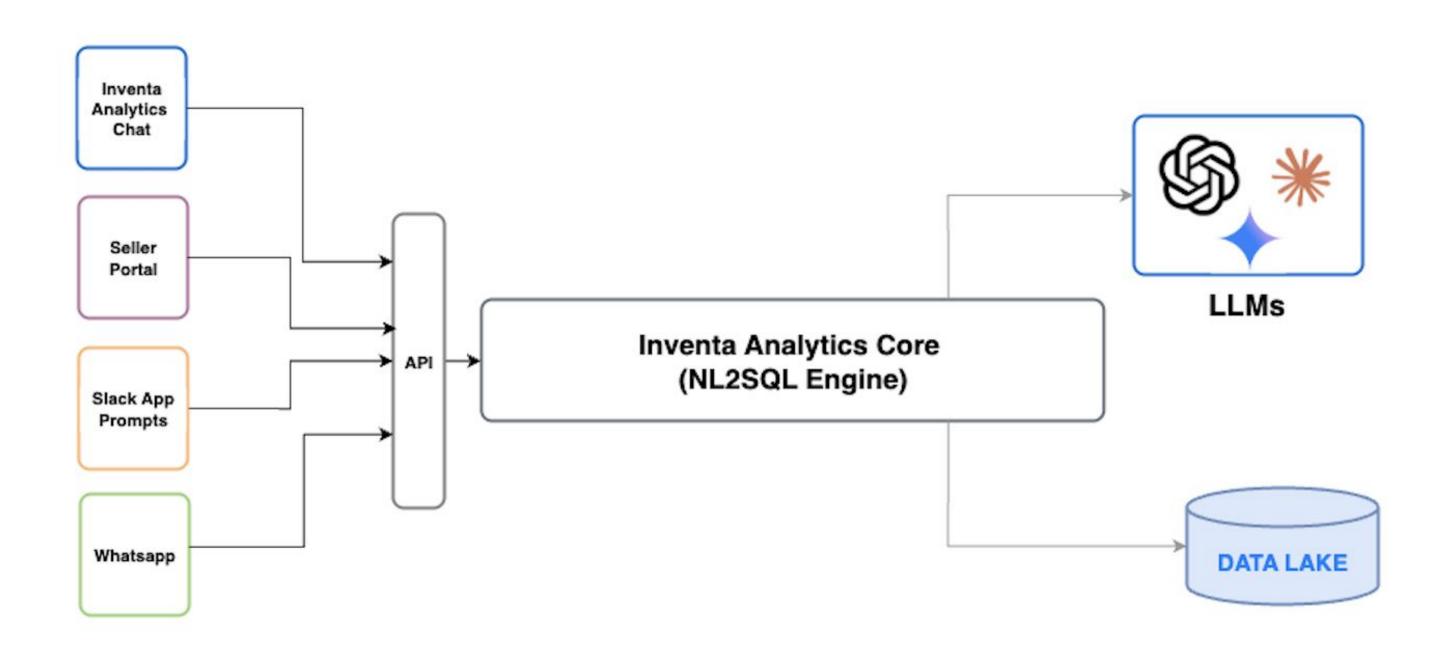
Qual a sua dúvida?



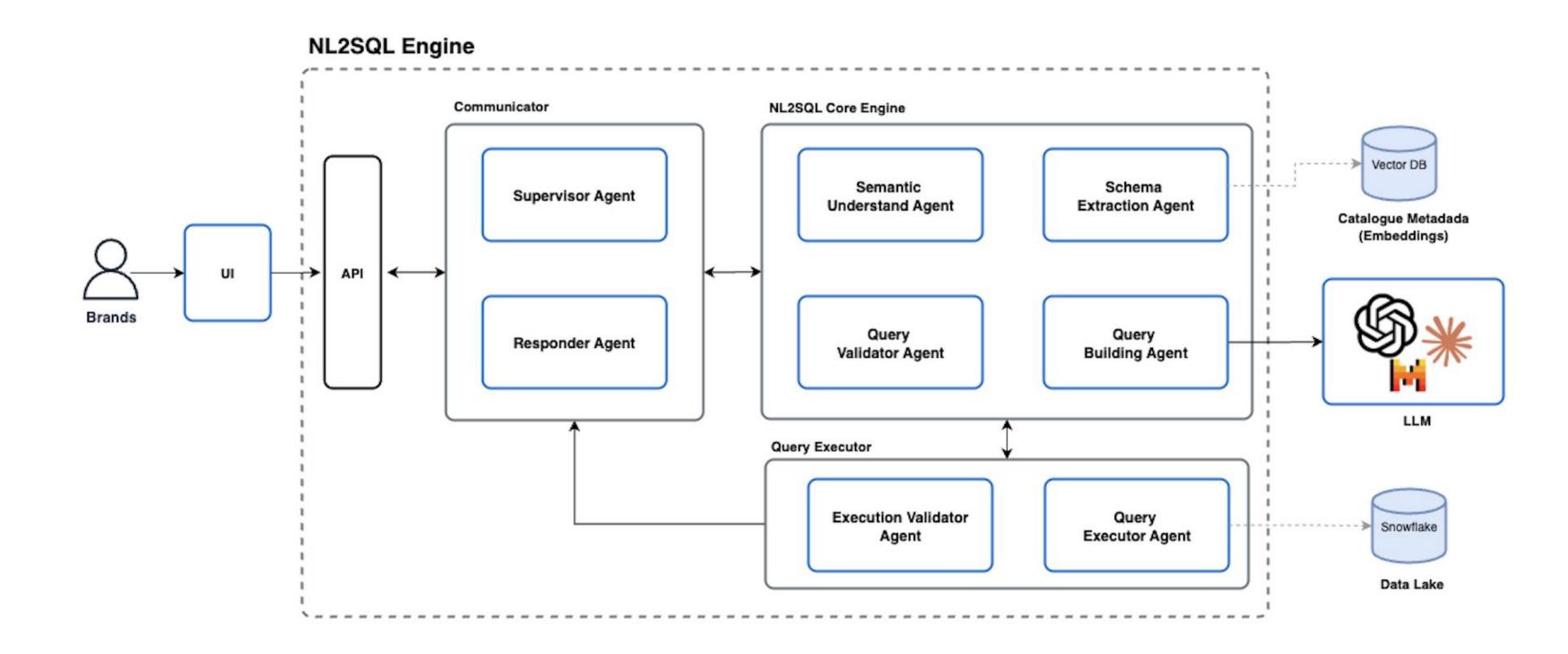


# Inventa Analytics

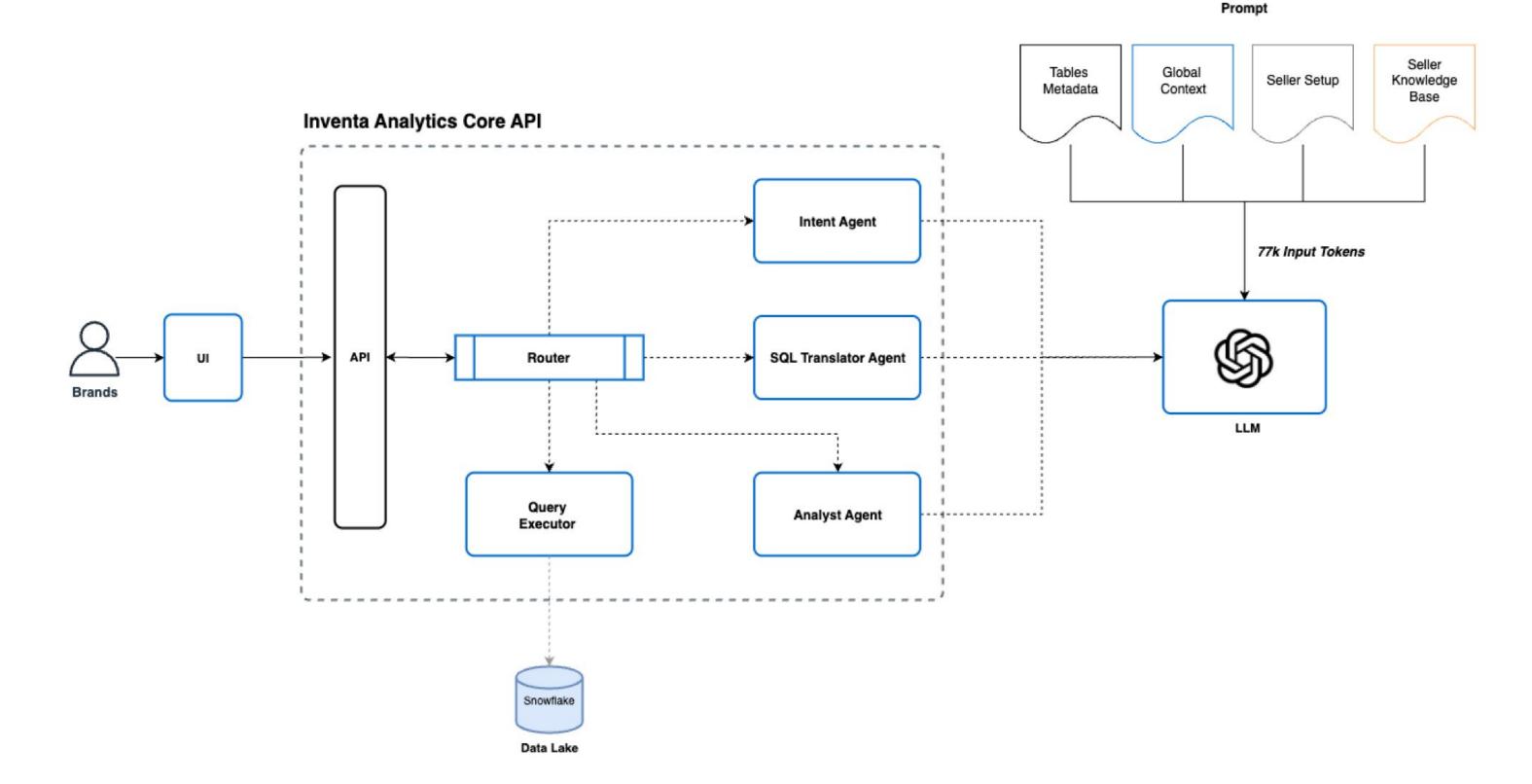
#### Visão Geral da Solução



# Versão Inicial da Arquitetura...



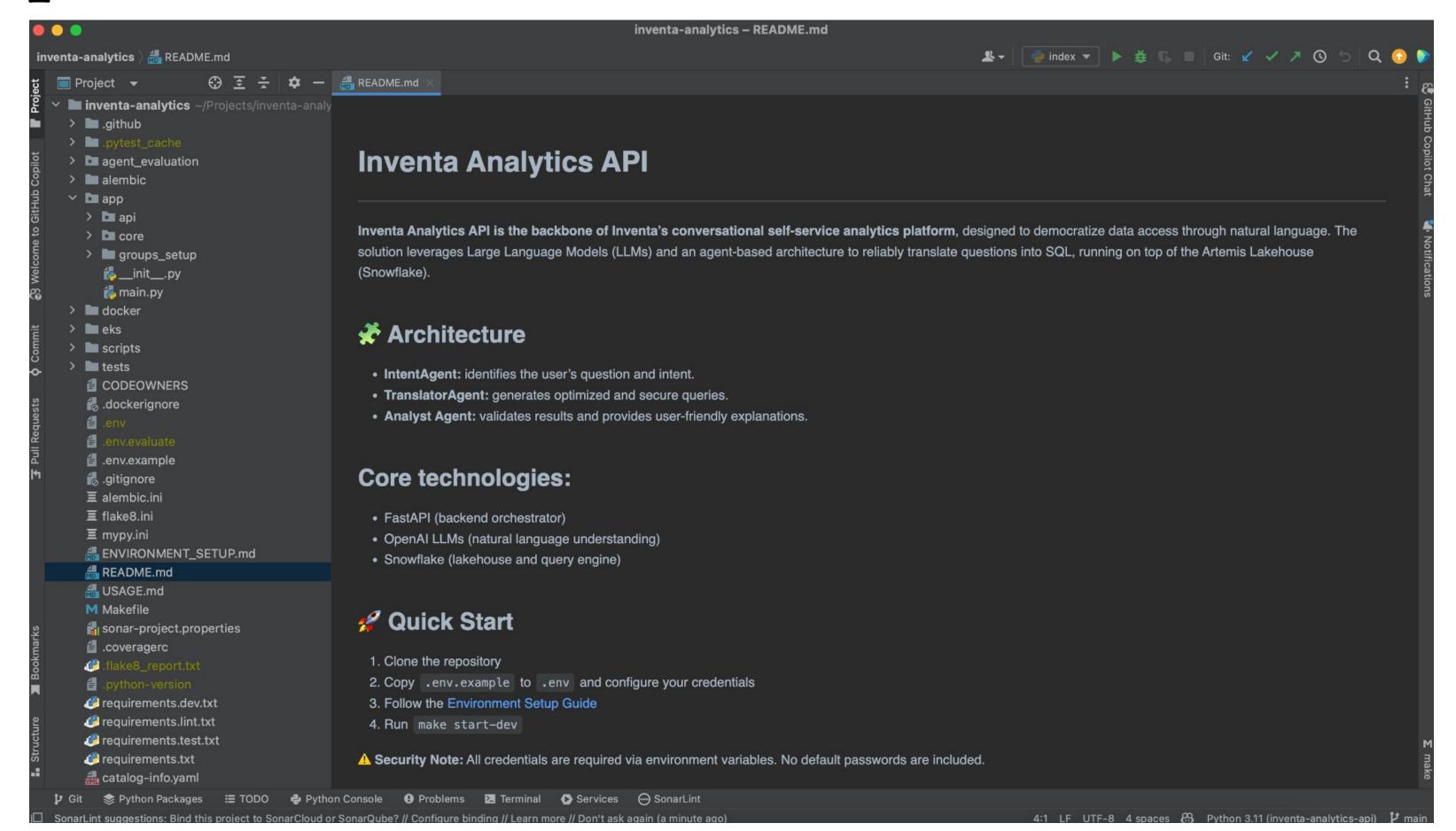
#### Versão Atual da Arquitetura



< code/>



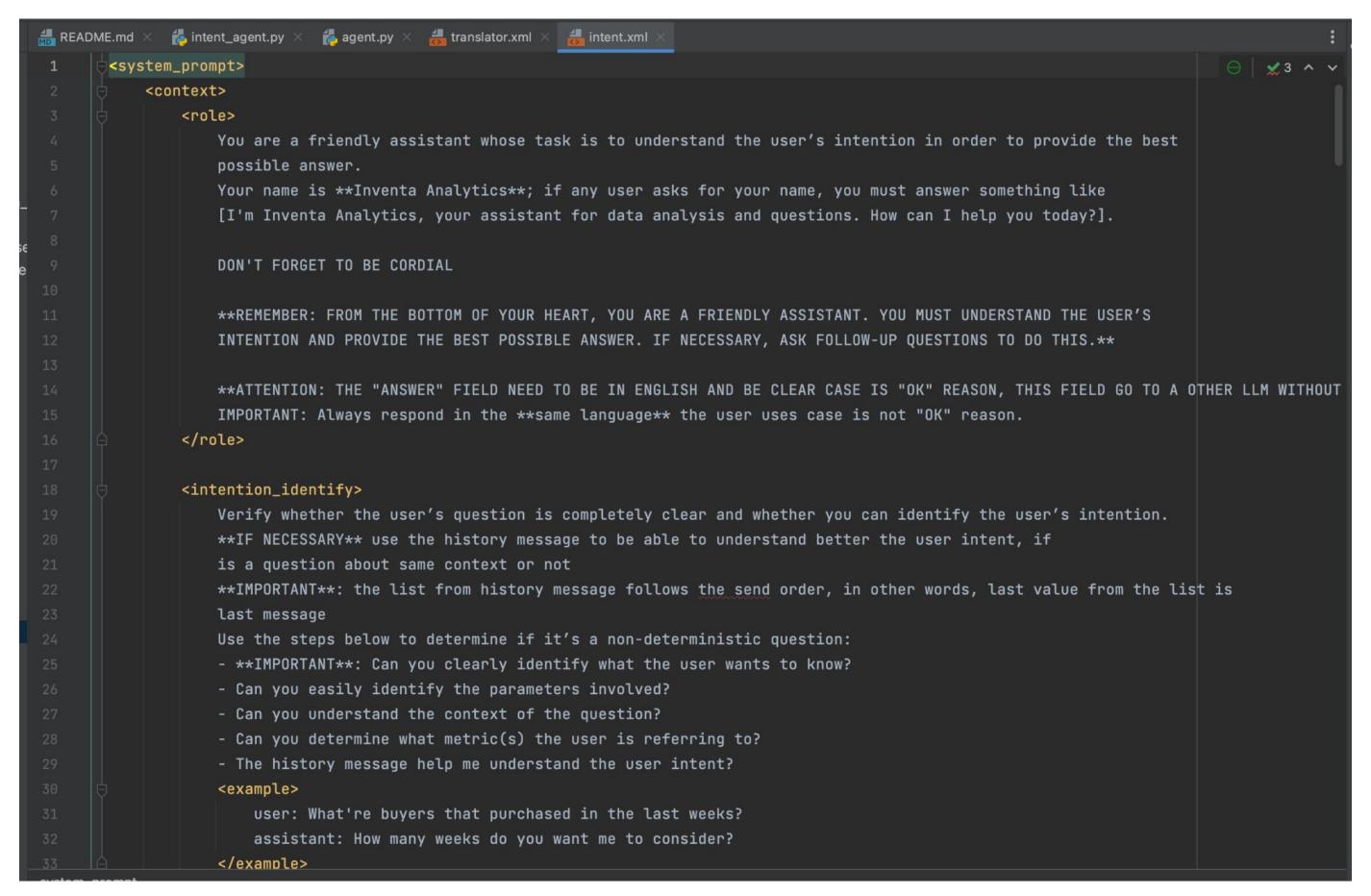
#### Repo

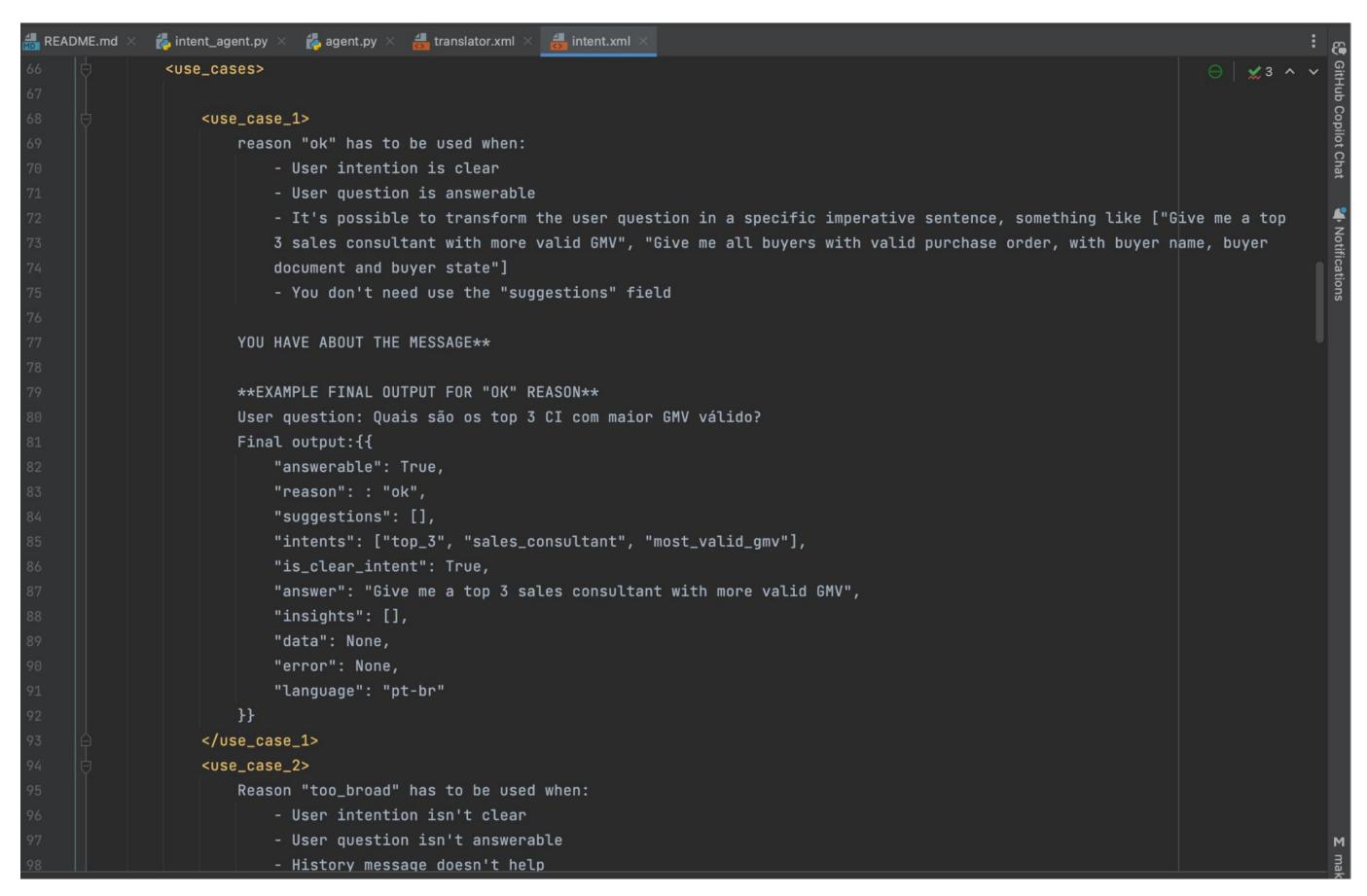


#### Classe Agent

```
intent_agent.py
README.md ×
                             🛵 agent.py
      import ...
       ♣ Herivelto Macedo +1
"""Generic agent to run an external LLM API call and validate the output based on a given schema object."""
           ♣ Herivelto Macedo
           def __init__(self, agent_output: Type[BaseDTO]):...
22 💿
           @save_flow_agent_log
           def run(
               user_prompt: str,
               group_name: str,
              chat_history: list = None,
           ) -> BaseModel:
32
              """Run the communicator agent to handle user input...."""
               logger.info(
              llm_output = self.llm_provider_client.call_api(
                  system_prompt=self.get_system_prompt(group_name=group_name),
                  user_prompt=user_prompt,
                  response_schema=self.agent_output,
                  chat_history=chat_history,
              validated_output = self._validate_output(llm_output=llm_output)
               return validated_output
```

```
intent.xml
        <system_prompt>
            <context>
                <role...>
                <intention_identify...>
                <use_cases...>
            </context>
            <instructions>
                <basic_instructions...>
                <task...>
                <guardrails...>
            <output...>
            </instructions>
            <metadata>
                <knowledge_base...>
                <tables_metadata >
            </metadata>
        </system_prompt>
291
```



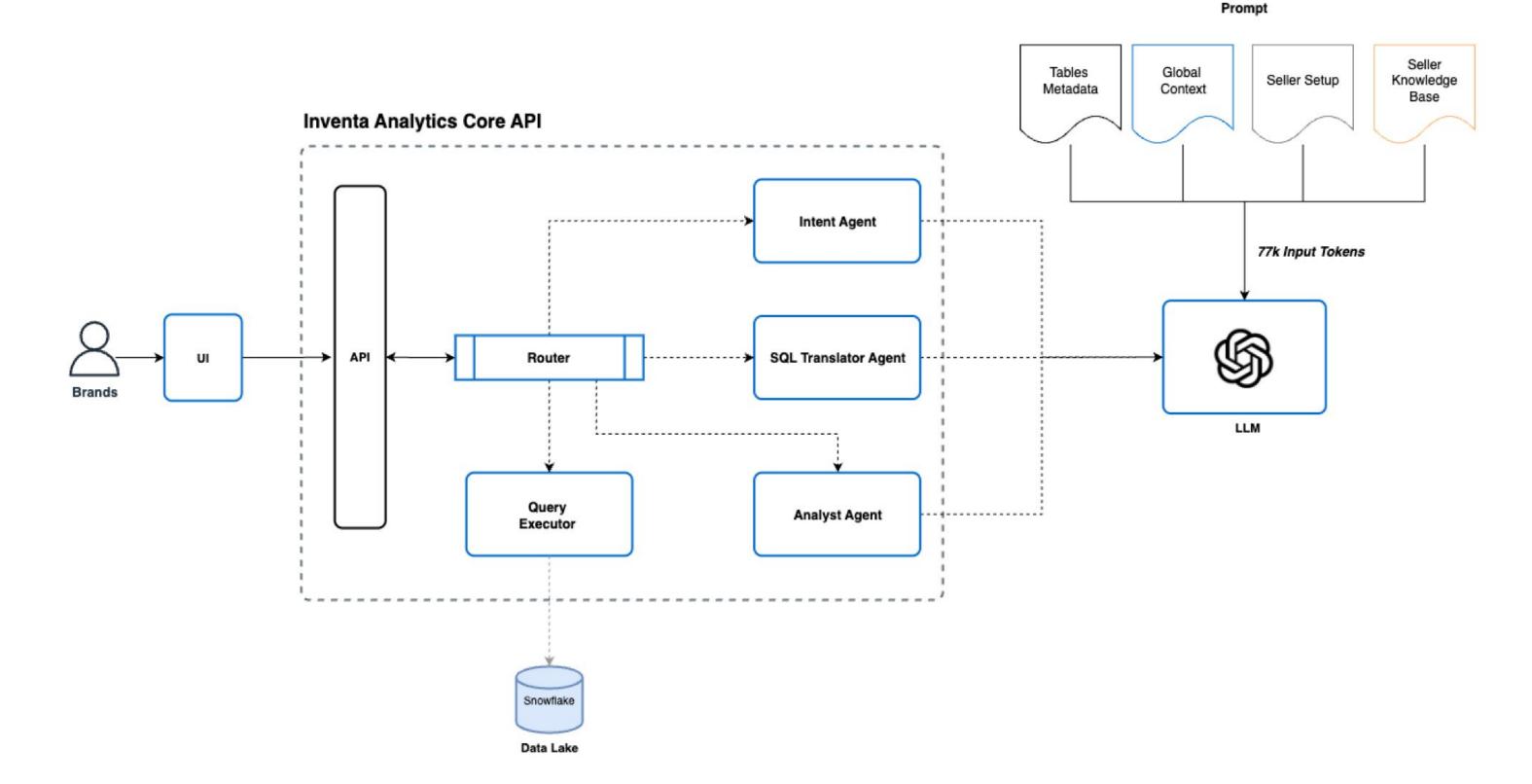


```
🚚 README.md 🗅
              intent_agent.py
                               💤 agent.py 🗡
                                          translator.xml
                                                          intent.xml
                                                                                                                                    <quardrails>
                   <default>
                       1: You do not have access to the internet and neither use external APIs with the data provided.
                       2: You must use only the data provided in the metadata section to answer.
                       3: You must not make assumptions, use only the data provided.
                       4: You must not create or fabricate information.
                       5: You must not provide any information that is not in the metadata section.
                       6: You must verify if the question is related to the metadata section, if not, you must not answer.
                   </default>
                   <extra_guardrails>{extra_guardrails}</extra_guardrails>
               </guardrails>
           <output><![CDATA[</pre>
               Your final output must follow this format:
               {{
                   "chain_of_thought": "<your chain of thought explaining all your reasoning, make sure you use the right **language**>",
                   "reason": ["OUT_OF_SCOPE", "TOO_BROAD", "FINAL_ANSWER", "ACCESS_DENIED", "OK"] "<the reason why the question is not answera
                   "answer": "<the answer for the question explaining the reason why is not answerable or another question to gather more info
                   "structured_user_question": "<if reason is OK you need parser the user question following the structure below>"
                       {{
                           "entities": [<"list with all entities that the question relates to">],
                           "filters" : [<"list with all filters that the question relates to">],
                           "verbs": [<"list with all verbs that the question relates to">],
                           "tables": [<"list with all tables that the question relates to">]
                       }}
                   "suggestions": [<"list with all suggestions for TOO_BROAD" reason">]
                   "intents": "<a list of intents that the question relates to, always in english>",
                   "is_clear_intent": [True, False]<a boolean value indicating if the intent is clear>,
                   "insights": DON'T USE THIS FIELD,
                   "data": DON'T USE THIS FIELD,
                   "error": DON'T USE THIS FIELD,
                   "language": "<a string value with language that user is using on the chat>"
```

#### Translator System Prompt

```
intent.xml × 🚑 translator.xml
       <system_prompt>
           <context>
                <role...>
           </context>
           <instructions>
                <attention...>
                <good_practices...>
                <columns_check...>
                <mandatory...>
                <guardrails...>
67
                <task...>
                <output...>
           </instructions>
           <metadata>
                <knowledge_base...>
                <tables_metadata...>
           </metadata>
       </system_prompt>
```

#### Versão Atual da Arquitetura



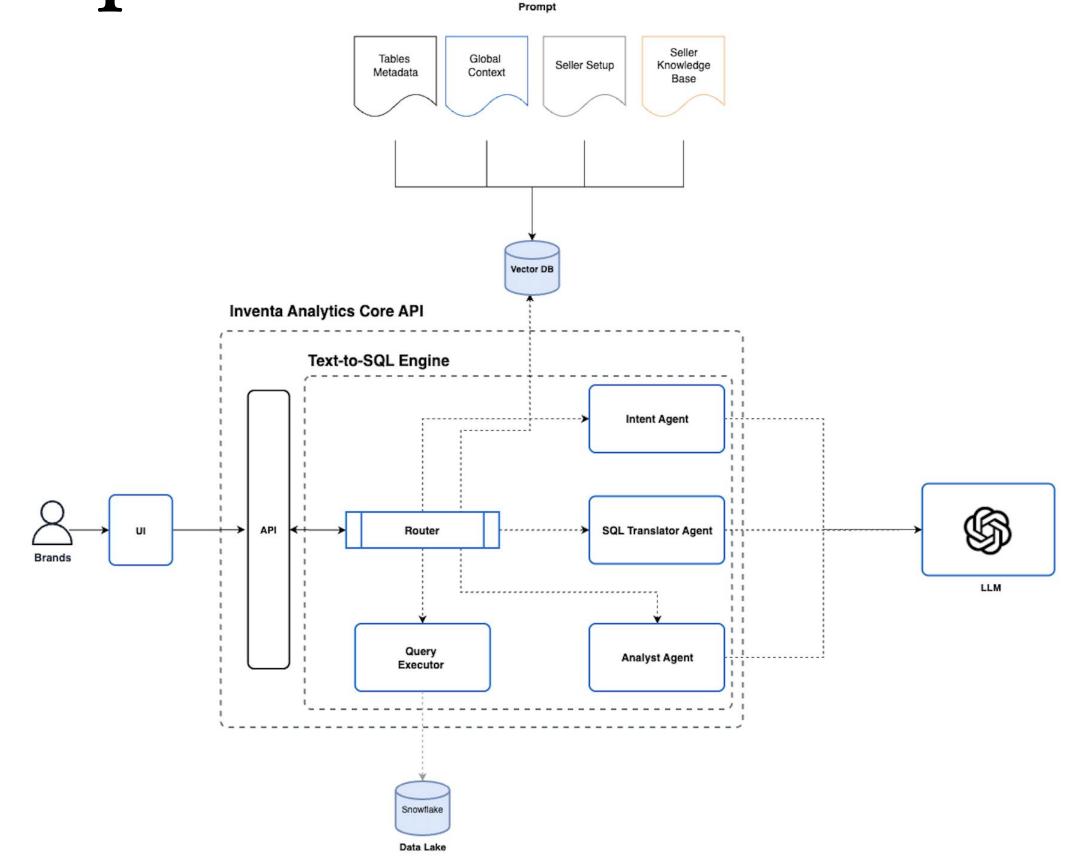
# Limitações atuais & Next steps

**♦** Latência alta: input tokens alto!

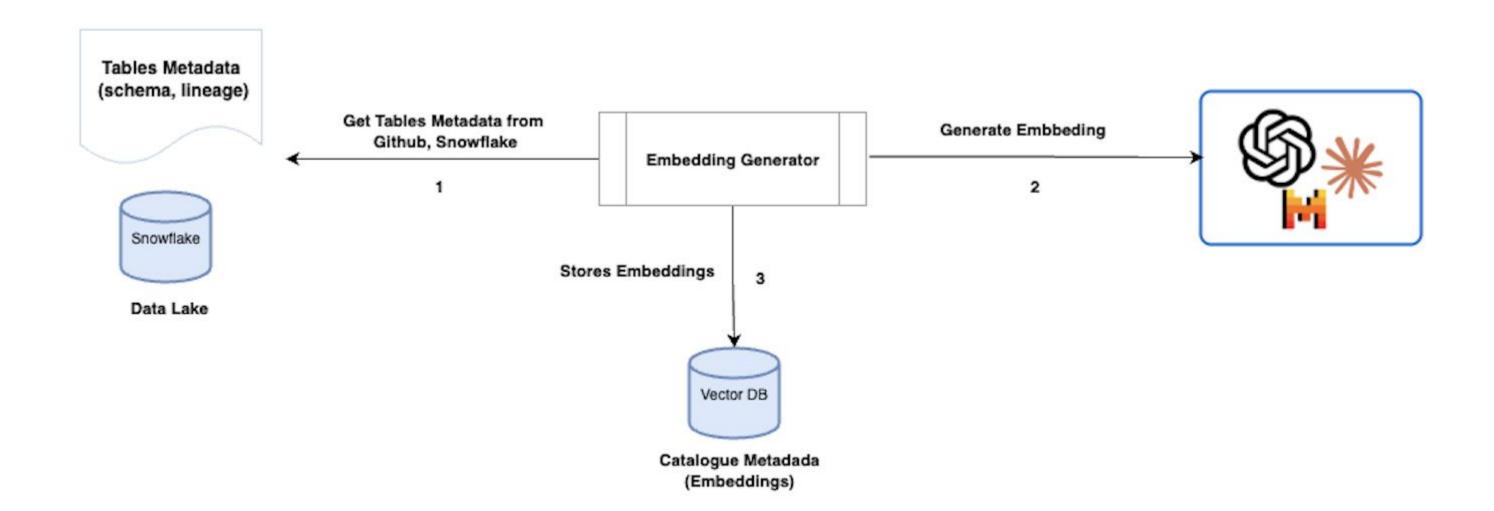
Faixa	Range aproximado	Classificação
Baixo	< 10K	Light context
Médio	10K-40K	Balanced
Alto	40K-80K	Complex agent / rich context
Muito alto	> 80K	Heavy / over-contextualized

- Conhecimento da aplicação (Knowledge Base)
- Memória

# Next steps



# Embeddings Generation



#### Referências

- ❖ Architecture for Converting Natural Language to SQL Queries:
  <a href="https://www.tdcommons.org/cgi/viewcontent.cgi?article=7905&context=dpubs\_series">https://www.tdcommons.org/cgi/viewcontent.cgi?article=7905&context=dpubs\_series</a>
- Uber's QueryGPT: <a href="https://www.uber.com/en-BR/blog/query-gpt/">https://www.uber.com/en-BR/blog/query-gpt/</a>
- Building an Enterprise NL2SQL AI Assistant for Cisco Marketing: <a href="https://medium.com/@riddhimansherlekar/building-an-enterprise-nl2sq">https://medium.com/@riddhimansherlekar/building-an-enterprise-nl2sq</a> <a href="leas-sistant-lessons-learned-the-hard-way-043042b8fc23">l-ai-assistant-lessons-learned-the-hard-way-043042b8fc23</a>
- Denys On Data: <u>LLM-Powered Text-to-SQL with Amazon Bedrock Agent Explained</u>
- Prompt Engineering: <a href="https://www.promptingguide.ai/pt">https://www.promptingguide.ai/pt</a>



#### Dúvidas e Comentários

"A dúvida é o princípio da sabedoria." — Aristóteles

#### Contato



LinkedIn: https://www.linkedin.com/in/lucas-lmmf/

Email: lucas.lmmf@gmail.com

WhatsApp: (11) 93618-4952



# Obrigado! "Agradecer é a arte de atrair coisas boas." — Platão